

## Mining Frequent Patterns, Associations, & Correlations

- **Motivation:** Business transaction records

- Discovery of *interesting correlation relationships* that help *business decision-making processes* (catalog design, cross-marketing, customer shopping behavior analysis, ...)

اللي خلانا نعمل فكرة mining frequent هو السوبر ماركت في الحالات العادية الزبون بيدخل السوبر ماركت ويشتري الحاجات اللي عايزها ويبدفع تمنها ويمشي وملهوش علاقي باللي بيحصل بعد كده المفروض بقي صاحب السوبر ماركت يشوف ايه اكثر حاجات اتباعت ويحطهم جنب بعض زي الشاي والسكر او العيش واللبن ويبدأ يعيد توزيعهم جنب بعد ولو لقي حاجه معينه محدش بيشتريها كثير ممكن يحطها جنبهم ويحدد يوم ايه اللي بيحصل فيه كده زي الاجازات كده ينزل الخصومات دي في يوم الاجازة بحيث يسهل عليه البيع

Cross marketing احط مثلا الموبايل جنب الاكسيسورات بتاعته عشان ازود فرصه شرائهم سوا

### Market basket analysis

مين اشترى ايه واشترى معاها ايه

- How to place SW, HW, and Accessories?

ودي علي حسب الشخص اللي بيبيع والمبيعات عنده يعني لو لقي HW , accessories بيتباعوا كثير ممكن يحط جنبهم SW عشان يزود نسبه بيعها يعني مسائله تقديرية

- Frequent patterns are item sets that appear frequently in a dataset (e.g. transaction records)

مجموعه items اللي بتكرر مع بعضها كثير

*Frequently associated* ممكن اطلع منها association rules يعني لو اشتريت كمبيوتر يبقى لازم اشترى انتي فيرس (اينعم احنا بنجيبه مسروق من النت بس ده في الدول المتقدمه)

- Support and Confidence are measures of rule interestingness

Support and Confidence بتحددلي مدي ثقتي في rule اللي انا حددتها دي

- 2% support → 2% of transactions show that computers and AV\_SW are bought together يشوف الحاجتين دول اتباعوا سوا
- 60% confidence → 60% of customers who bought a computer also bought AV\_SW ان الزبون اللي اشترى كمبيوتر اشترى معاها انتي فيرس

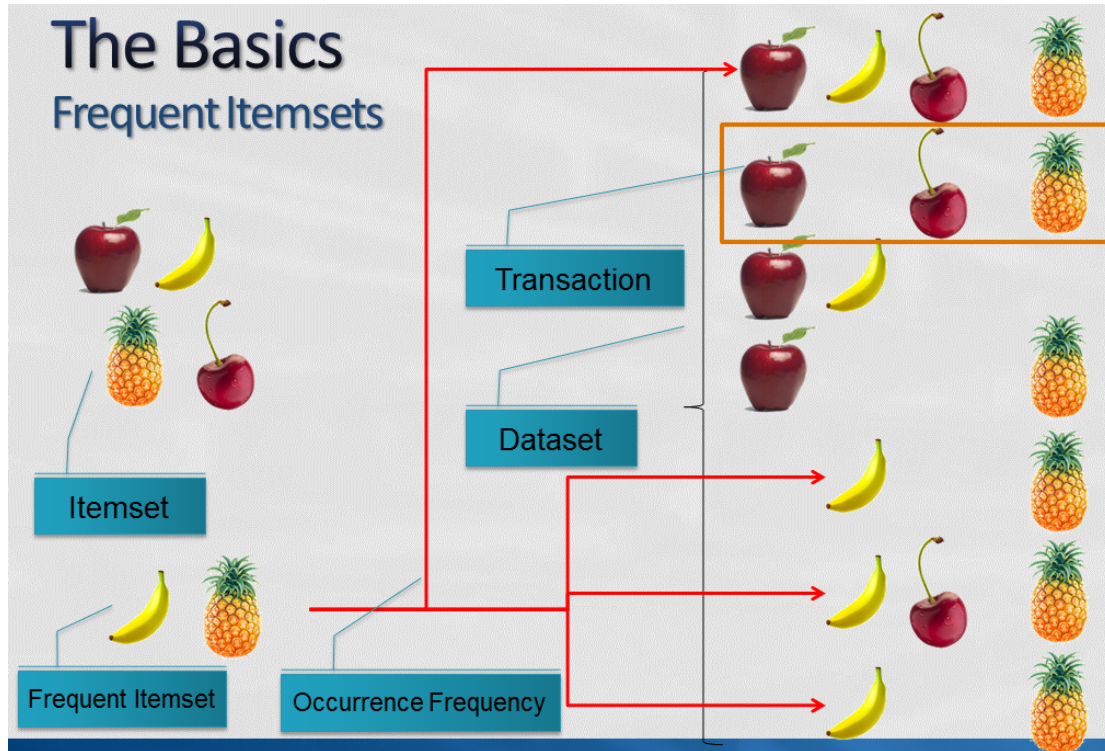
---

كل صف من الصورة بيمثل transaction

كل الانواع الموجوده هي itemsets

كل الصفوف مع بعض بيمثل Dataset

عدد مرات ظهور الحاجات اللي بيدور عليها مع بعض بتسمي Occurrence Frequency



عدد تكرار العنصرين مع بعض أكثر من threshold معين هو frequent itemset يعني الموز والananas اتكرروا كام مرة لو انا عامل threshold=2 يبقى عشان اقول عليهم frequent itemset لازم يكونوا اتكرروا مرتين او أكثر

$$\text{support}(A \Rightarrow B) = P(A \cup B) = \frac{n(A \cup B)}{N}$$

Support هو عدد المرات اللي اتكرر فيها item معينه علي عدد itemsets كلها يعني هنا عدد مرات تكرار الموز والananas علي عدد transaction كلها هابقي  $4/7=57\%$

$$\text{support}(\text{banana} \Rightarrow \text{pineapple}) = P(\text{banana} \cup \text{pineapple}) = \frac{4}{7} = 57\%$$

$$\text{confidence}(A \Rightarrow B) = P(B|A) = \frac{n(A \cup B)}{n(A)}$$

Confidence هو كام واحد من اللي اشتري item الاول اشتري item الثاني يعني كام واحد اشتري موز اشتري معاه اناناس هنلاقي عدد النا اللي اشترت الاتنين مع بعض 4 وعدد الناس اللي اشترت الموز 5

$$\text{confidence}(\text{banana} \Rightarrow \text{pineapple}) = P(\text{pineapple}|\text{banana}) = \frac{4}{5} = 80\%$$

min\_support count اقل عدد اتكرر عشان اقول عليه itemset

min\_support and min\_confidence انا اللي بحددهم علي حسب احتياجي

🔹 If a rule satisfies min\_support and min\_confidence thresholds, it is said to be **strong**

لو rule اللي كنت حددتها قبل كده اكبر من min\_support and min\_confidence thresholds او بتساويهم بيتقال عليها strong

● problem of mining association rules reduced to mining frequent itemsets

● Association rules mining becomes a two-step process:

1. Find all **frequent itemsets** with frequently  $\geq$  a predetermined *min\_support count*   
 ودي تعتبر اكثر خطوه مكلفه
2. Generate **strong association rules** from the frequent itemsets that satisfy *min\_support* and *min\_confidence*

● If *min\_support* count is set too low  $\rightarrow$  huge # of frequent itemsets

لو رقم *min\_support* count قليل هايبيقي عندي عدد كبير جدا من frequent itemsets

## Mining Frequent Itemsets

### 1) Apriori Algorithm

بعتد علي حاجه بتكون كوجوده مسبقا قبل ماطبق الالجوريزم

لو عندي item مش frequent لو زودت عليه item ثاني مش هايغير من حالته بمعني

A C

A C D

B

لو threshold=2 بيقي B مش frequent فلو زودت اي حرف ثاني جنب B وليكن C مش هايخليها frequent

١. هابص علي كل item لوحد اتكرر كام مرة

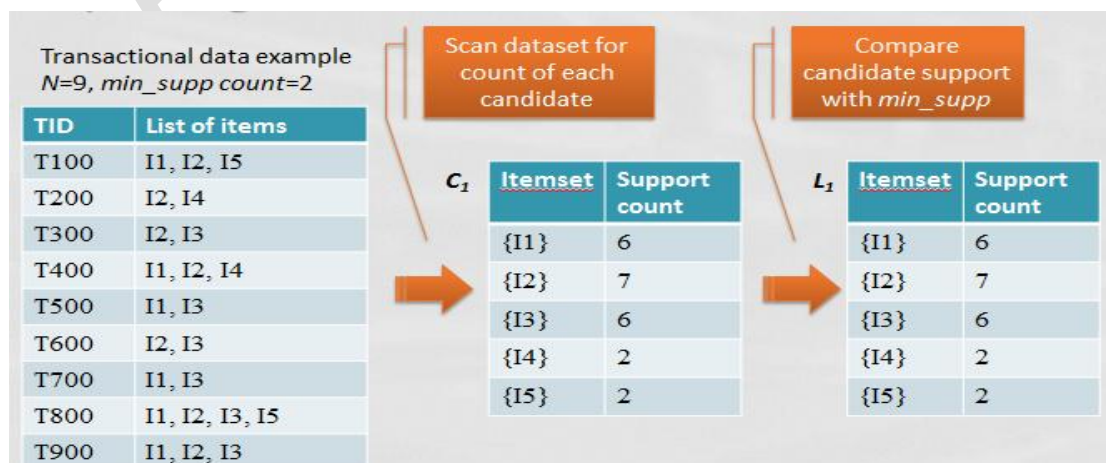
٢. هعمل generate لحاجه اسمها candidate itemset

To improve efficiency, use the **Apriori property**:

“All nonempty subsets of a frequent itemset must also be frequent” – if a set cannot pass a test, all of its supersets will fail the same test as well

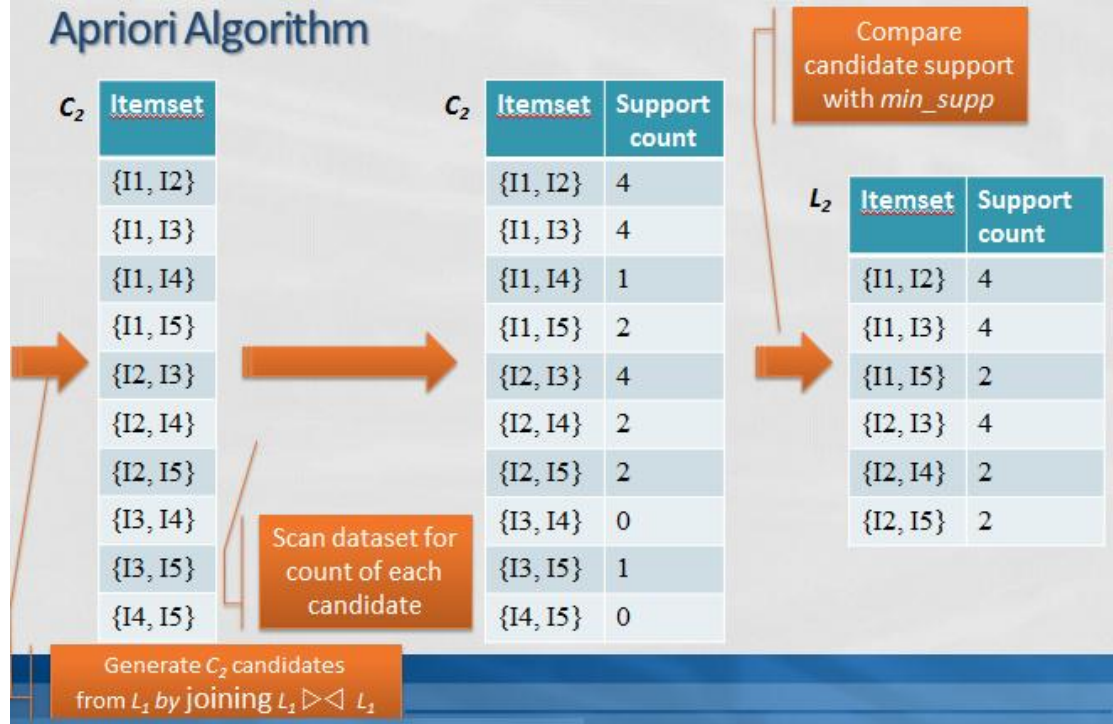
if  $P(I) < \text{min\_support}$  then  $P(I \cup A) < \text{min\_support}$

#### Level 1



١. اول جدول هو كل Transactions اللي عندي
٢. ثاني جدول بشوف كل item اكرر كام مرة يعني I1 اكررت ٦ مرات في T100 , T400 , T500 , T700 , T800 , T900
٣. ثالث جدول اي Transactions اقل من threshold=2 اشيله هلاقي هنا كلهم اكبر من او يساوي ٢

## Level 2

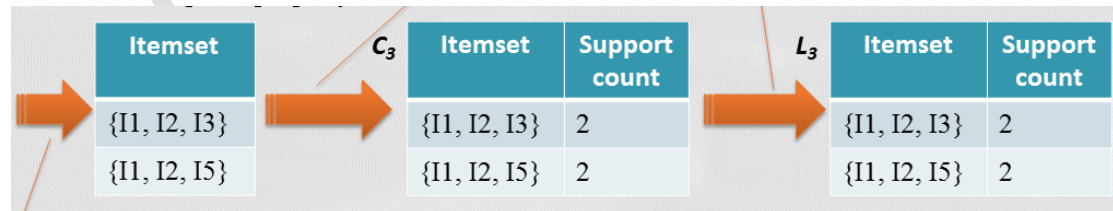


بيدا اخذ كل اثنين مع بعض واشوفهم اكررنا مع بعض:-

١. اول جدول بعمل تبادل وتوافق لكل item مع الباقي
٢. ثاني جدول بشوف كل اثنين item اكررنا قد ايه في الجدول الاصلي اللي فيه transactions
٣. بعد كده بشيل اي itemset اكررنا اقل من ٢

## Level 3

باخذ كل صفين من الجدول اللي في ليفل ٢ واعملهم join بس لازم يكونوا مشتركين في اول item يعني اخذ I1,I2 مع I1,I3 فهما مشتركين في I1 وهو اول item فيهم بيقى اعملهم join وارتيهم في العموم لما بختار كل join بيبقي لازم يكونوا مشتركين في كل items ماعدا الاخير



(١) بعد ما بختار كل صفين مشتركين في اول item هلاقيهم

{I1, I2, I3}, {I1, I2, I5}, {I1, I3, I5}, {I2, I3, I4}, {I2, I3, I5}, {I2, I4, I5}

بس مش هأخذ منهم الا {I1, I2, I3}, {I1, I2, I5} بس لان جميع عناصرهم frequent يعني I1,I2 و I1,I3 و I2,I3 والتاني {I1, I2, I5} هلاقي I1,I2 و I1,I5 و I2,I5 هما اللي frequent (اللي support count اكبر من 2

Not all subsets are frequent هاحذفهم بسبب {I1, I3, I5}, {I2, I3, I4}, {I2, I3, I5}, {I2, I4, I5} اول واحد  
مثلا {I1, I3, I5} frequent بس I3,I5 مش frequent <= I1,I5 و I1,I3 frequent وهكذا في الباقي

#### Level 4



مش هلاقي حاجه frequent فاهقف علي اخر حاجه كنت فيها frequent في level 3

Itemset
{I1, I2, I3}
{I1, I2, I5}

### Generating Association Rules from Frequent Itemsets

Association rules can be generated using the *confidence equation*, as follows:

عشان احسب association rule من frequent item بحسبها بالقانون ده

$$\text{confidence}(A \Rightarrow B) = P(B|A) = \frac{\text{support\_count}(A \cup B)}{\text{support\_count}(A)}$$

- For each frequent itemset  $l$ , generate all nonempty subsets of  $l$
- For every nonempty subset  $s$  of  $l$ , output the rule:

$$s \Rightarrow (l - s) \quad \text{if} \quad \frac{\text{support\_count}(l)}{\text{support\_count}(s)} \geq \text{min\_confidence}$$

هاخدكل صف في اخر جدول الجدول مثلا الصف الثاني وهي {I1,I2,I5} هابدا اقسامها ل subset واطلع منها  
association rule هأخذ اول عنصرين {I1,I2}بيشاوروا علي العنصر الثالث I5 وهكذا زي الرسمه

Nonempty subsets	Association Rules
{I1, I2}	{I1, I2} $\Rightarrow$ I5
{I1, I5}	{I1, I5} $\Rightarrow$ I2
{I2, I5}	{I2, I5} $\Rightarrow$ I1
{I1}	{I1} $\Rightarrow$ {I2, I5}
{I2}	{I2} $\Rightarrow$ {I1, I5}
{I5}	{I5} $\Rightarrow$ {I1, I2}



حاسب **confidence** بالقانون اللي فوق اللي هو  $\text{confidence}(A \Rightarrow B) = P(B|A) = \frac{\text{support\_count}(A \cup B)}{\text{support\_count}(A)}$

في البسط عدد ظهور كل العناصر مع بعض وفي المقام عدد ظهور nonempty subsets يعني علي اول واحده عدد ظهور {I1,I2,I5} علي عدد ظهور {I1,I2} في الجدول الاصلي اللي في الشمال تحت عدد مرات ظهور العناصر كلها 2= اللي هما T100 و T800 وعدد مرات ظهور I1,I2 = 4 اللي هما T100 , T400 , T800 , T900

$$\text{confidence}(A \Rightarrow B) = P(B|A) = \frac{2}{4} = 50\%$$

Itemset
{I1, I2, I3}
{I1, I2, I5}

وبعد ما احسب confidence اشوف **min\_confidence** وهكذا لباقي العناصر في الجدول واحذف اي confidence اقل منه

**min\_confidence** = 70%

TID	List of items	Nonempty subsets	Association Rules	Confidence
T100	I1, I2, I5			
T200	I2, I4	{I1, I2}	{I1, I2} ⇒ I5	2/4 = 50%
T300	I2, I3	{I1, I5}	{I1, I5} ⇒ I2	2/2 = 100%
T400	I1, I2, I4	{I2, I5}	{I2, I5} ⇒ I1	2/2 = 100%
T500	I1, I3	{I1}	{I1} ⇒ {I2, I5}	2/6 = 33%
T600	I2, I3	{I2}	{I2} ⇒ {I1, I5}	2/7 = 29%
T700	I1, I3	{I5}	{I5} ⇒ {I1, I2}	2/2 = 100%
T800	I1, I2, I3, I5			
T900	I1, I2, I3			

ودول اللي هأخدهم من الجدول هنا لان confidence اكبر من 70%

{I1, I5} ⇒ I2	2/2 = 100%
{I2, I5} ⇒ I1	2/2 = 100%
{I5} ⇒ {I1, I2}	2/2 = 100%

الاجوريزم ده شغال كويس بس لو عندي item كتير هاتبقى مشكله فهاستخدم الجوريزم FP-Growth

## 2) FP-Growth

- To avoid costly candidate generation

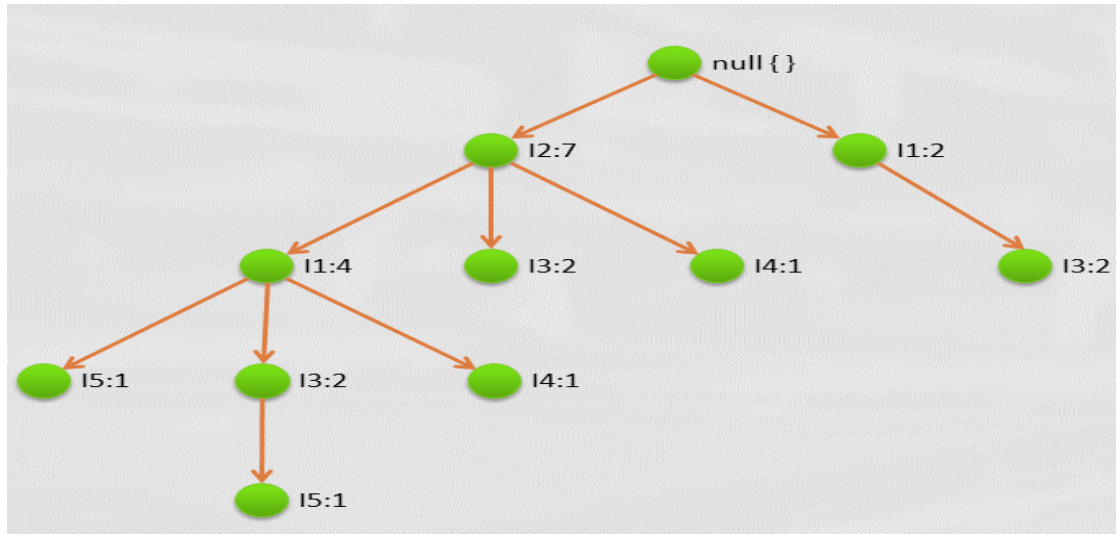
وهاستخدمه عشان مش كل شوية اعمل candidate

- Divide-and-conquer strategy:

1. Compress database representing frequent items into a **frequent pattern tree (FP-tree)** – 2 passes over dataset
2. Divide compressed database (FP-tree) into *conditional databases*, then mine each for frequent itemsets – traverse through the FP-tree



وهكذا لحد ماخلص كل الجدول وفي الاخر هاتكون tree كده

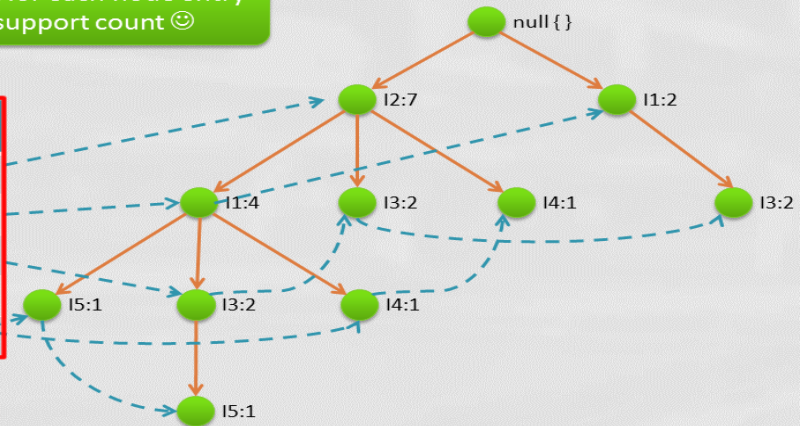


Trace the node link path for each node entry and you get that item's support count 😊

**L<sub>1</sub> - Reordered**

Itemset	Support count	Node Link
{I2}	7	
{I1}	6	
{I3}	6	
{I4}	2	
{I5}	2	

For Tree Traversal



كل item في tree هاشاور عليه في node link يعني I2 موجوده مرة واحده ممكن تكون اكرر اكثر من مرة بس مرسومه مرة واحد اشاور من node link علي النقطة دي I1 اترسمت مرتين يبقي اطلع سهم من node بتاع I1 بيمر بكل نقط I1 وهكذا في كلهم

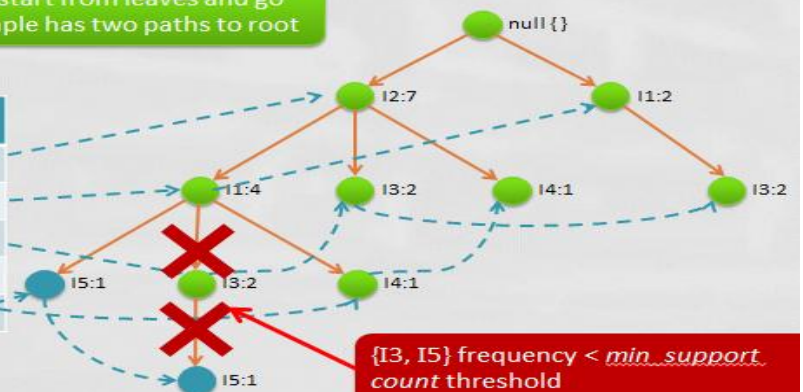
الاجوريزم ده Bottom-up algorithm يعني هابدا باخر نقطه تحت خالص اللي هي I5 لحد لما اوصل لل root

I5 هاشوف frequency هلاقيها اقل من  $\text{min\_support count threshold}$  فهاشيله

Bottom-up algorithm – start from leaves and go up to root – I5 for example has two paths to root

**L<sub>1</sub> - Reordered**

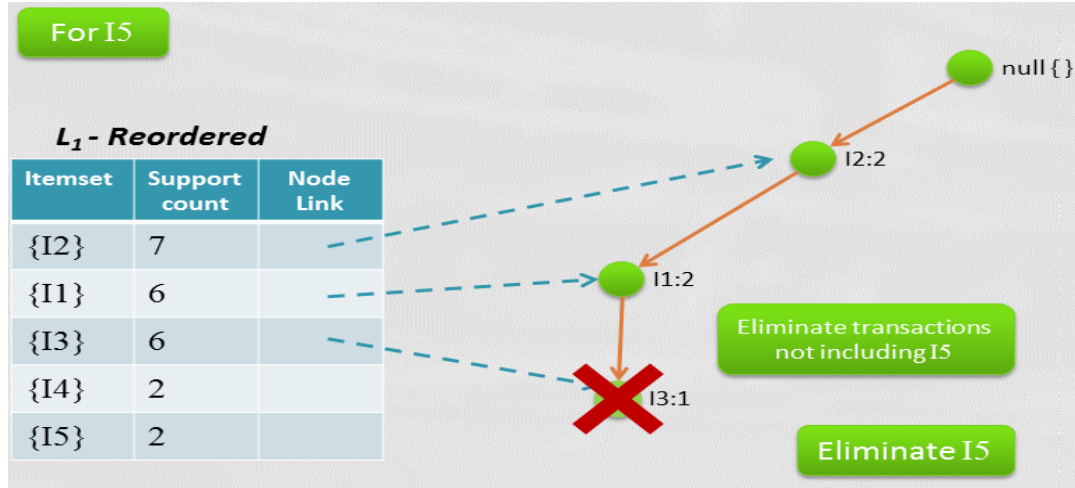
Itemset	Support count	Node Link
{I2}	7	
{I1}	6	
{I3}	6	
{I4}	2	
{I5}	2	





هاعمل Conditional FP-tree Construction لكل node لوحدها

هاشيل كل transaction اللي مفيهاش node دي ولو تكن اول حاجه i5 هايتبقي T800 , T100 هاشيل منهم I5 وارسم اللي باقي فهاقي بعد مامسحت I5 (I1,I2) T100 و (I1,I2,I3) T800 هارسمهم مع بعض زي ماعملت فوق وامسح اي node اتكررت اقل من min support



I2 اكثر item اتكرر مرتين مرة مع T100 و مرة مع T800 وكذلك I1 بس I2 في الترتيب الاول كانت رقم واحد فهانا هاتبقي الاول وبعد كده I1 وبعدهم I3 اللي اتكررت مرة واحده بس في T800 فهانشيلها لانها اقل من min support

1) Condition pattern base

وده الاحتمالات الممكنة قبل مالمسح I3

2) Condition FP-tree

وده node اللي موجوده بعد الحذف وعدد تكرار كل واحده

3) Frequent patterns generated

ودي المسارات الممكنة عشان اروح ل I5 بعد ماشيلت I3

وامسك كل item واعمل عليه نفس العملية دي وهايطلع الجدول ده

Item	Conditional Pattern Base	Conditional FP-tree	Frequent Patterns Generated
I5	{ {I2, I1: 1}, {I2, I1, I3: 1} }	$\langle I2: 2, I1: 2 \rangle$	{I2, I5: 2}, {I1, I5: 2}, {I2, I1, I5: 2}
I4	{ {I2, I1: 1}, {I2: 1} }	$\langle I2: 2 \rangle$	{I2, I4: 2}
I3	{ {I2, I1: 2}, {I2: 1}, {I1: 2} }	$\langle I2: 4, I1: 2 \rangle, \langle I1: 2 \rangle$	{I2, I3: 4}, {I1, I3: 4}, {I2, I1, I3: 2}
I1	{ {I2: 4} }	$\langle I2: 4 \rangle$	{I2, I1: 4}

Paths to which item is suffix      Prefix paths to item after eliminating infrequent items

## Pattern Evaluation Methods


- Not all association rules are interesting

مش شرط كل association rule اللي طلعتها تبقي مفيدة ومهمه ممكن مايكونش ليها لازمه

- buys(X, "computer games")  $\Rightarrow$  buys(X, "videos") [40%, 66%]
- P("videos") is already 75% > 66%
- The two items are negatively associated  $\rightarrow$  buying one *decreases* the likelihood of buying the other

- We need to measure "*real strength*" of rule

- Correlation analysis**

Correlation  الاتنين item بيزيدوا مع بعض او بيقلوا مع بعض اوا ملهمش علاقه ببعض

$$A \Rightarrow B [\text{support, confidence, correlation}]$$

عشان اعمل evaluation لل association rule اللي طلعتها بازود parameter كمان وهو correlation بلاضافه الي support و confidence

$$1. \text{ lift} = \frac{P(A \cup B)}{P(A)P(B)}$$

- A and B are independent if  $P(A \cup B) = P(A)P(B)$

لو  $P(A \cup B) = P(A)P(B)$  بيقى قسمتهم بتساوي ١ ويبقى A, B مش معتمدين علي بعض

- Otherwise, dependent and correlated occurrence

لو مش بتساوي واحد بيقى معتمدين علي بعض وفيه correlation مابينهم

- If  $\text{lift} < 1$ , A is negatively correlated with B

لو الناتج اصغر من ١ بيقى A negatively correlated with B

- If  $\text{lift} > 1$ , A is positively correlated with B  $\rightarrow$  A's occurrence "**lifts**" the occurrence of B

لو اكبر من ١ بيقى A positively correlated with B  $\rightarrow$  A's occurrence "**lifts**" the occurrence of B

- 2.  $\chi^2 \rightarrow$  already discussed in a previous lecture

تاني حاجه chi square ودي خدناها في قبل كده